**SAN DIEGO STATE UNIVERSITY**

**Antoni Luque Santolaria <aluque@sdsu.edu>**

# More "super-short" Caudovirales data

**Antoni Luque Santolaria** <aluque@sdsu.edu>          Thu, Jan 14, 2021 at 12:58 PM
To: Stephen Nayfach <snayfach@lbl.gov>

Thanks, Stephen.

1. Let me see if I understood it. The length in the biomes that you
sent me is what I would obtain if I filter the dataset DTRs_20kb.csv
for each biome-category and sum the contig lengths for each case. Is
that right?

2. The length of the DTRs in the original file that you sent me ranged
from 0.11% to 17% of the contig lengths. The average is 1.24%, and the
median is 0.90%. Adding the DTR length, thus, would not impact the
current genome length analysis, but I'll add a comment in the draft so
we can review this point to be accurate. It would be interesting,
however, to analyze the direct terminal repeats of the final
candidates. But we can do that in the next round.

3. When I analyzed the dataset, I didn't dereplicate the contigs (for
some reason I thought the data was dereplicated). When scrutinizing
the proteins of the 11 potential candidates for T=1 capsids (viral
circular contigs with MCP and < 5 kbp), I found that three of the
contigs are (quasi) identical:  DTR_770701, DTR_758375, and
DTR_742353. The first two are 100% identical. The third one differs in
two bases. I think that having assemblies of three (quasi) identical
contigs is a good sign. The three came from the human gut. Do you know
if they're from the same sample?

4. The analysis of the nanopore datasets sounds really exciting. We
can coordinate on that after I finish the draft that I'm preparing.

Best,
Toni

On Wed, Jan 13, 2021 at 12:31 PM Stephen Nayfach <snayfach@lbl.gov> wrote:
>
> Hi Antoni,
>
>> Got it, Stephen. Thanks for the quick response. I'll let you and Simon
>> judge if a map would be compelling when you get a chance to evaluate
>> the draft that I'm preparing. If needed, the coordinates associated

>> with the contigs that came from the Global Ocean Virome should be
>> accessible. I assume that combining the IMG and GOV would be a good
>> chunk of original data.
>
>
> Sounds like a good plan.
>>
>>
>> Here is another request. Let's see if this one is more doable. Do you
>> have easy access to a variable, for example, genome content (in
>> bases), to compare the contribution of each biome-level (1, 2, 3) to
>> the complete dataset where you extracted the circular genomes (<20
>> kb)?
>
>
> Unfortunately, not. But I do have access to the full set of circular genomes (>2kb) for each biome, which I could use to calculate the cumulative genome length per biome. Would that be helpful?
>
>>
>> I want to test if the frequency of circular contigs (<20 kb) observed
>> across biomes is similar to the proportion of the genetic content that
>> each biome contributed to the total dataset. The null hypothesis is
>> that both distributions would be similar. The assumption behind the
>> hypothesis is that all biomes would have a similar relative content of
>> small phages. Significant differences between the distributions of
>> small circular contigs vs. the genetic contribution of a biome to the
>> total dataset would indicate that some biomes are richer (or poorer)
>> in circular small phages.
>
>
> We can certainly test the ratio of small vs large Caudovirales phages using the information on hand.
>
> Also, thinking back to a previous comment you made, all the sequences I sent did have terminal repeats (I believe I clipped the end off in the final dataset) and should certainly be viral. So we may want to go back and include some of those that have been excluded. But, I'm also happy to proceed with the current set you have for the initial analysis, as we can always go back later to add more sequences. For example, I've now assembled close to 4,000 nanopore datasets, which could be really useful to give us more confidence these small genomes are bona fide.
>
>> Total genomic content: 30,000 Mb (Megabases)
>>
>> Biome level 1:
>> Host-associated --> 25,000 Mb
>> Aquatic --> 10,000 Mb
>> ...
>>
>> Let me know if this makes any sense to you and if it is possible.
>
>
> Please see attached for numbers of total genome length of circular Caudovirales from each biome.
>

> Best,
> Stephen


\-\-
Dr. Antoni Luque
http://luquelab.com
Google Scholar profile
Assistant Professor
Department of Mathematics and Statistics
Viral Information Institute
Computational Science Research Center
San Diego State University